

Fernández JM, de-la-Torre V, Richardson D, Royo R, Puiggròs M, Moncunill V, Fragkogianni S, Clarke L, BLUEPRINT Consortium, Flicek P, Rico D, Torrents D, Carrillo-de-Santa-Pau E, Valencia A. [The BLUEPRINT Data Analysis Portal](#). *Cell Systems* 2016, 3(5), 491-495.e5.

Copyright:

© 2016 Elsevier Inc. Published with open access.

DOI link to article:

<http://dx.doi.org/10.1016/j.cels.2016.10.021>

Date deposited:

06/01/2017

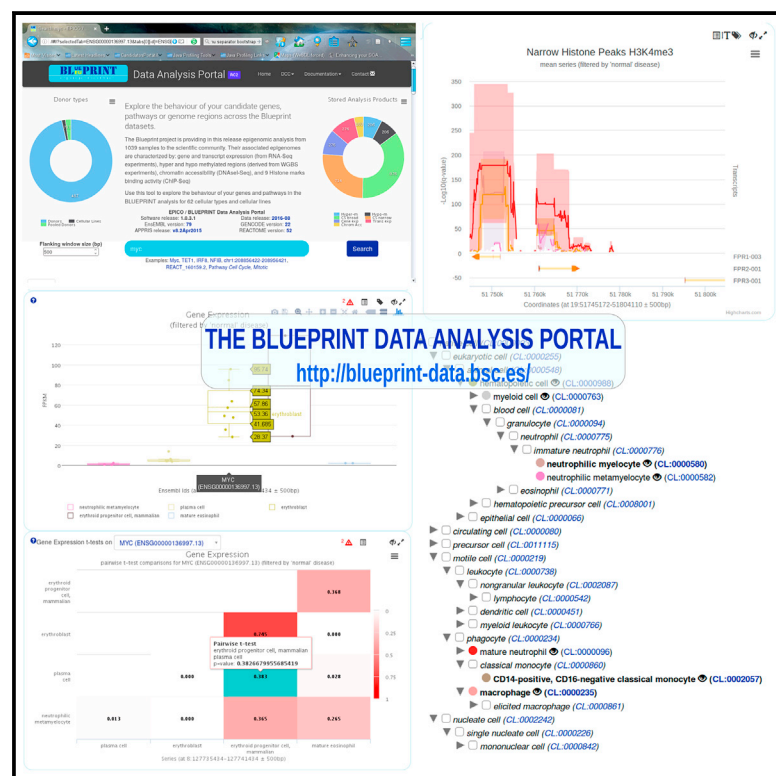


This work is licensed under a [Creative Commons Attribution 4.0 International License](#)

Cell Systems

The BLUEPRINT Data Analysis Portal

Graphical Abstract



Authors

José María Fernández,
Victor de la Torre, David Richardson, ...,
David Torrents,
Enrique Carrillo de Santa Pau,
Alfonso Valencia

Correspondence

ecarrillo@cnio.es (E.C.d.S.P.),
valencia@cnio.es (A.V.)

In Brief

As part of the International Human Epigenome Consortium, this study introduces EPICO, an epigenomics comparative cyber-infrastructure to develop comparative epigenomic web portals to display large epigenomic datasets. EPICO is used to implement the BLUEPRINT Data Analysis Portal (BDAP) for user-friendly access to BLUEPRINT results. Explore the Cell Press IHEC web portal at <http://www.cell.com/consortium/IHEC>.

Highlights

- EPICO infrastructure generates data-mining portals for large epigenomic series
- EPICO was used to implement the user-friendly BLUEPRINT Data Analysis Portal (BDAP)
- Users can retrieve epigenomic processed data demonstrated with FPR1 and IRF8
- Data and plots can be easily downloaded for downstream analysis and publication



The BLUEPRINT Data Analysis Portal

José María Fernández,^{1,2} Víctor de la Torre,^{1,2} David Richardson,³ Romina Royo,^{2,4} Montserrat Puiggròs,⁴ Valentí Moncunill,⁴ Stamatina Fragkogianni,⁴ Laura Clarke,³ BLUEPRINT Consortium,⁶ Paul Flicek,³ Daniel Rico,^{1,7} David Torrents,^{4,5} Enrique Carrillo de Santa Pau,^{1,*} and Alfonso Valencia^{1,2,8,*}

¹Structural Biology and BioComputing Programme, Spanish National Cancer Research Centre (CNIO), Madrid 28029, Spain

²Spanish Bioinformatics Institute INB-ISCIII ES-ELIXIR, Madrid 28029, Spain

³European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

⁴Barcelona Supercomputing Center (BSC), Joint BSC-CRG-IRB, Research Program in Computational Biology, BSC - CRG - IRB, Barcelona 08028, Spain

⁵Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona 08010, Spain

⁶<http://www.blueprint-epigenome.eu/>

⁷Present address: Institute of Cellular Medicine, Newcastle University, Newcastle upon Tyne NE2 4HH, UK

⁸Lead Contact

*Correspondence: ecarrillo@cnio.es (E.C.d.S.P.), valencia@cnio.es (A.V.)

<http://dx.doi.org/10.1016/j.cels.2016.10.021>

SUMMARY

The impact of large and complex epigenomic datasets on biological insights or clinical applications is limited by the lack of accessibility by easy, intuitive, and fast tools. Here, we describe an epigenomics comparative cyber-infrastructure (EPICO), an open-access reference set of libraries to develop comparative epigenomic data portals. Using EPICO, large epigenome projects can make available their rich datasets to the community without requiring specific technical skills. As a first instance of EPICO, we implemented the BLUEPRINT Data Analysis Portal (BDAP). BDAP provides a desktop for the comparative analysis of epigenomes of hematopoietic cell types based on results, such as the position of epigenetic features, from basic analysis pipelines. The BDAP interface facilitates interactive exploration of genomic regions, genes, and pathways in the context of differentiation of hematopoietic lineages. This work represents initial steps toward broadly accessible integrative analysis of epigenomic data across international consortia. EPICO can be accessed at <https://github.com/inab>, and BDAP can be accessed at <http://blueprint-data.bsc.es>.

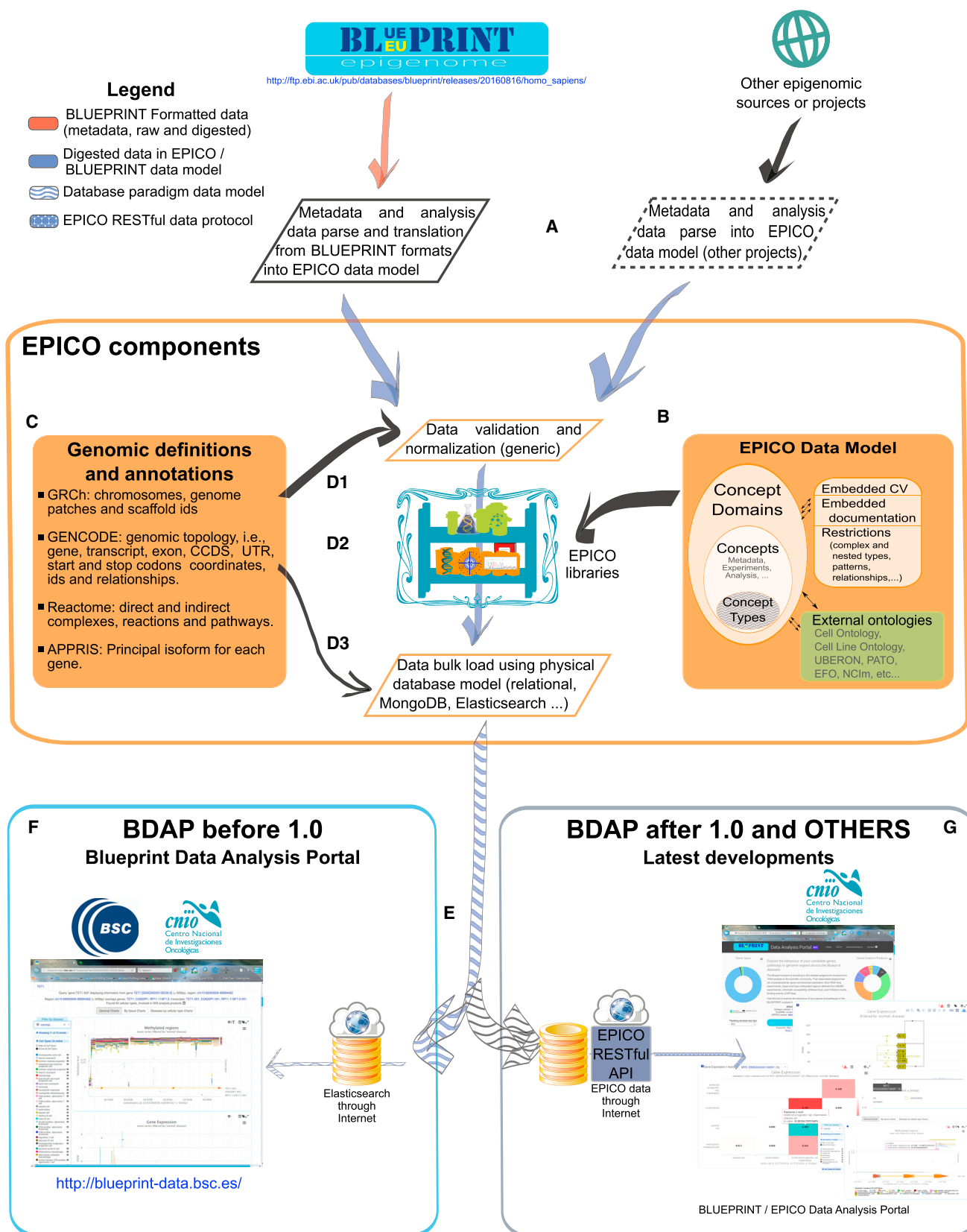
The International Human Epigenome Consortium (IHEC; [IHEC, 2016](#)) coordinates standards for the production, distribution, and accessibility of reference epigenomes generated by several large consortia, including BLUEPRINT ([Adams et al., 2012](#)), CEMT ([CEMT, 2016](#)), CREST ([CREST, 2016](#)), DEEP ([DEEP, 2016](#)), ENCODE ([ENCODE, 2016](#)), CEEHRC ([CEEHRC, 2016](#)), and NIH ROADMAP ([Roadmap Epigenomics Project, 2016](#)). Each consortium keeps its original data on their own Data Coordination Center (DCC) portal, which provides some additional analysis results (for example, chromatin state or intron retention) in different formats (text, BED, BigWig for the raw signal, and

BigBed for regions highly enriched in raw signal) along with metadata ([Table S1](#)). Moreover, some consortia provide genome browsers to visualize and compare the primary data ([Table S1](#)). Nevertheless, additional bioinformatics skills are needed to identify, download, process, and analyze the large and complex datasets ([Table S1](#)). Indeed, the majority of potential users interested in epigenomic datasets, including most biologists and physicians, are not able to exploit the data satisfactorily. Therefore, novel efficient exploration tools to quickly test biological hypotheses are needed.

Here, we describe an epigenomics comparative cyber-infrastructure (EPICO; <https://github.com/inab>) to facilitate the production of user-friendly interactive web portals, and we describe a portal for BLUEPRINT data (<http://blueprint-data.bsc.es>) implemented using EPICO. The BLUEPRINT Consortium is a flagship European project that aims to provide reference epigenomes from hematopoietic cell lineages ([Adams et al., 2012](#)). Portals created with EPICO enable the comparison of the epigenetic structure of different cell types and related diseases.

The EPICO platform is based on five components: (1) a data model ([EPICO-data-model, 2016](#)); (2) data validation and loading programs ([EPICO-data-loading-scripts, 2016](#)), which must be adapted to the data and metadata acquisition of the particular project; (3) an empty database that will store all the data and metadata produced by the data validation and loading programs; (4) the EPICO REST API ([EPICO-REST-API, 2016](#)), which implements the queries to the database, providing a programmatic access; and (5) the data analysis portal itself ([BP-Data_analysis, 2016](#)), which queries the databases through the EPICO REST API ([Figure 1](#)).

The minimum infrastructure needed to generate epigenomic data portals with EPICO are the five aforementioned components, storage space to create the database, a connection to fetch the primary data to be integrated into the database, and the modules to receive the queries over the stored data and send the results to be visualized. A detailed technical description of the EPICO components and the instructions to create custom data portals are provided at <https://github.com/inab/epico-data-analysis-portal/wiki>.



(legend on next page)

The epigenomic information displayed by BDAP is based on data obtained from chromatin immunoprecipitation sequencing (ChIP-seq) experiments (H3K4me3, H3K4me1, H3K27ac, H3K36me3, H3K27me3, H3K9me3, or H2A.Zac); DNase-seq (DNase I sequencing); WGBS (whole-genome bisulfite sequencing of hypo- and hyper-methylated regions); and RNA-seq (RNA sequencing; at the gene or transcript level). As of the 2016-08 data release, the platform contains the analysis of 2,757 products from 2,558 experiments, involving 487 donors, 11 pool donors, and seven cell lines, and summarizing a total of 62 different cell types that cover 17 diseases. The BDAP allows users to visualize and compare all the epigenomic and transcriptomic data for blood cell types of interest. This query is performed following the three-step process implemented in EPICO (STAR Methods). We illustrate its utility by analyzing two genes as examples, Formyl Peptide Receptor 1 (*FPR1*; Figures S1 and S2) and interferon regulatory factor 8 (*IRF8*; Figures S3 and S4).

FPR1 is one of the most extensively studied G-protein-coupled receptors involved in neutrophil chemotaxis (reviewed by Ye et al., 2009). We used the BDAP to explore the phase in which neutrophil differentiation *FPR1* expression is regulated (STAR Methods). A gene expression boxplot shows a clear increase in *FPR1* expression, as neutrophil differentiation progresses from the neutrophilic myelocyte to the mature neutrophil, and a more modest increase for *FPR2* (Figures S2A and S2B). The cell types with the strongest expression were the segmented neutrophils of the bone marrow and the mature neutrophils. In addition, the ChIP-seq data revealed active histone modifications in the start codon and transcribed regions of the principal isoforms for *FPR1* and *FPR2* in the cell types with the strongest expression (Figures S2C and S2D). These results suggest that the chemotactic properties associated with *FPR1* and *FPR2* are acquired during neutrophil differentiation and that they reach their peak in the segmented neutrophils of the bone marrow.

The interferon regulatory factor 8 (*IRF8*) has been identified as a key transcription factor that regulates myeloid cell production. *IRF8* maintains the balance between monocytes and neutrophils, and a lack of this gene increases the number of neutrophils and diminishes the monocyte population (Kurotaki et al., 2014). We used the BDAP to explore the differences between neutrophils and monocytes at the transcriptome and epigenome levels (STAR Methods). Gene expression boxplots clearly show that *IRF8* is expressed in classical monocytes and macrophages,

whereas neutrophils do not express this transcription factor (Figure S4B). In addition, differences in the transcripts expressed are observed (Figure S4A). Moreover, peaks of active histone modifications H3K27ac were observed along the gene body and at the start codon in the two cell types that express *IRF8* (Figure S4D), while in neutrophilic myelocytes, repressive histone modifications were situated in the region around the start codon (H3K27me3 and H3K9me3; Figures S4C and S4E). These results confirmed previous observations suggesting that *IRF8* is a marker of the macrophage lineage (Kurotaki et al., 2014).

Data portals facilitate access to different data analyses by reducing the need to deal with issues related to the different formats in which data are stored. Unfortunately, the current lack of standards for data and metadata in epigenomics limits the possibility of developing a single portal to compare data of the different IHEC consortia. Given this situation, we propose the use of EPICO to create project-specific data analysis portals. EPICO includes a common template and standards for the description of data and metadata.

The BLUEPRINT Data Analysis Portal complements the BLUEPRINT-DCC portal (DCC_portal, 2016) by providing the facilities to analyze multiple epigenetic data types at once—for instance, DNA methylation and histone marks—and to deal with multiple samples from different cell types rather than dealing with individual samples and specific data types. Moreover, the BDAP answers queries about specific genomic regions, genes, or pathways, providing summary statistics and comparative analysis grouping samples by cell type or tissue of origin. The BDAP is accessible to many biologists and doctors without programmatic or technical skills. This is different from other solutions, such as the one recently proposed by DeepBlue (Albrecht et al., 2016; Table S1). The BDAP is the first platform generated with EPICO. The main condition for generating a portal is that data and metadata have to be converted to EPICO standard format (STAR Methods). In the future, additional efforts will have to be made to homogenize across IHEC experimental procedures, quality controls, and primary data analysis workflows, a situation reminiscent of the current developments in other large-scale consortia, such as ICGC (ICGC, 2016) or ENCODE (ENCODE, 2016).

Future improvements to our platform will include the integration of our visualization tools (e.g., boxplots or scatterplots) and results (e.g., consensus peaks) with standard genomic browsers, such as Ensembl (Flicek et al., 2014) or the UCSC (University of California, Santa Cruz) Genome Browser (Kent et al., 2002).

Figure 1. EPICO Infrastructure Flowchart

- (A) Each epigenomic dataset usually has its own file formats and conventions, so this step is custom.
- (B) EPICO data model concepts, ontologies, and restrictions are common. Only details like the versions of reference Ensembl, GENCODE, GRCh, and other primary database resources or project name have to be tweaked.
- (C) As genomic definitions and annotations are published in common sites, and their data formats are stable from release to release, this step is done by EPICO.
- (D) The metadata and data insertion (which should be following the EPICO data model at this point) is composed of several steps, all of them generic: data validation and normalization (1) using EPICO libraries (2), which later translate it into the dependent database model (3) (currently supported relational, MongoDB, and Elasticsearch). In the case of BLUEPRINT, we have used Elasticsearch.
- (E) The data are massively inserted into the database, which already contains the database definitions mapped from the EPICO data model, as well as the ontologies, and the genomic coordinates of the known features, like genes, transcripts, direct complexes, reactions, and pathways.
- (F) The BLUEPRINT Data Analysis Portal prior to version 1.0 was issuing its queries to the read-only instance of Elasticsearch, which contained all the BLUEPRINT metadata + primary analysis data.
- (G) BDAP 1.0 issues its queries to the EPICO REST API, which manages the different databases and implements the queries to Elasticsearch. The EPICO Data Analysis Portal is going to be a superset of BDAP, able to work with one or more project datasets at once. Data from different epigenomic projects usually cannot be mixed and compared due to different experimental, normalization, and analysis protocols.

In summary, EPICO provides the infrastructure and a standard template to create powerful tools that bring complex epigenomic data to the hands of researchers who want to test biological hypotheses, as shown in the two use cases of the BDAP implementation of EPICO fed with BLUEPRINT primary data.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **CONTACT FOR REAGENT AND RESOURCE SHARING**
- **EXPERIMENTAL MODEL AND SUBJECT DETAILS**
- **METHOD DETAILS**
 - EPICO Cyber-Infrastructure Description
 - BLUEPRINT Data Analysis Portal General Usage Description
 - Step-by-Step Example of BDAP Usage: FPR1
 - Step-by-Step Example of BDAP Usage: IRF8
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
- **DATA AND SOFTWARE AVAILABILITY**
- **ADDITIONAL RESOURCES**

SUPPLEMENTAL INFORMATION

Supplemental Information includes four figures and one table and can be found with this article online at <http://dx.doi.org/10.1016/j.cels.2016.10.021>.

AUTHOR CONTRIBUTIONS

J.M.F. designed and wrote BP-Schema-tools, BP-DCC-model, EPICO-REST-API, and EPICO-data-loading-scripts and most of the BDAP web client, as well as the technical part of the manuscript. V.d.I.T. designed and wrote the first prototype of the BDAP web client. D. Richardson and L.C. provided insightful discussions and advice about the data model behind the BDAP, as well as feedback on the different BDAP releases. R.R., M.P., V.M., and S.F. have been responsible for the different public releases of the BDAP at BSC, synchronized with BLUEPRINT data releases, and they provided feedback on the data model. E.C.d.S.P. and D. Rico prepared the practical cases, beta tested the BDAP web client, and suggested improvements, as well as coordinating the preparation of the manuscript. E.C.d.S.P. coordinated the communication with the journal and the response to the reviewers, as well as the different manuscript versions. P.F., D.T., and A.V. coordinated the project.

ACKNOWLEDGMENTS

The authors thank David Pisano and Miriam Rubio from UBio-CNIO for their contributions at the early stages of the conceptual data model used by EPICO, as well as Avik Data (EMBL-EBI) for feedback on each BLUEPRINT data release. The INB-CNIO unit is a member of ProteoRed, PRB2-ISCIII, and is supported by grant PT13/0001, from the PE I+D+i 2013-2016, funded by ISCIII and FEDER. BSC-CRG-IRB acknowledges the funding support of the Spanish Ministry of Health, ISCIII, in the project Instituto Nacional de Bioinformática - PRB2: PT13/0001/0028. The research leading to these results was funded by the European Union's Seventh Framework Programme (FP7/2007-2013), under grant agreement number 282510 (BLUEPRINT), and by the European Molecular Biology Laboratory and the Spanish National Bioinformatics Institute.

Received: July 15, 2016

Revised: October 10, 2016

Accepted: October 24, 2016

Published: November 15, 2016

REFERENCES

- Adams, D., Altucci, L., Antonarakis, S.E., Ballesteros, J., Beck, S., Bird, A., Bock, C., Boehm, B., Campo, E., Caricasole, A., et al. (2012). BLUEPRINT to decode the epigenetic signature written in blood. *Nat. Biotechnol.* **30**, 224–226.
- Albrecht, F., List, M., Bock, C., and Lengauer, T. (2016). DeepBlue epigenomic data server: programmatic data retrieval and analysis of epigenome region sets. *Nucleic Acids Res.* **44** (W1), W581–W586.
- Bard, J., Rhee, S.Y., and Ashburner, M. (2005). An ontology for cell types. *Genome Biol.* **6**, R21.
- BP-analysis (2016). BLUEPRINT analysis descriptions release 20160816. ftp://ftp.ebi.ac.uk/pub/databases/blueprint/releases/20160816/homo_sapiens/.
- BP-Data_analysis (2016). BLUEPRINT data analysis portal GitHub repository. <https://github.com/inab/epico-data-analysis-portal>.
- BP-FPKM (2016). A description file about how FPKMs were calculated in release 20160816. ftp://ftp.ebi.ac.uk/pub/databases/blueprint/releases/20160816/homo_sapiens/README_maseq_analysis_crg_20160816.
- BP-Schema-tools (2015). Bioinformatic Pantry Schema tools GitHub repository. <https://github.com/inab/BP-Schema-tools>.
- CEEHRC (2016). McGill Epigenomics Mapping Centre. <http://epigenomesportal.ca/edcc/>.
- CEMT. (2016). Canadian Epigenetics, Environment and Health Research Consortium (CEEHRC) Network. <http://www.epigenomes.ca/>.
- Chen, P.P.-S. (1976). The entity-relationship model—toward a unified view of data. *ACM Trans. Database Syst.* **1**, 9–36.
- Codd, E.F. (1979). Extending the database relational model to capture more meaning. *ACM Trans. Database Syst.* **4**, 397–434.
- CREST (2016). International human epigenome consortium, IHEC, team Japan. <http://crest-ihec.jp/english/index.html>.
- Croft, D., Mundo, A.F., Haw, R., Milacic, M., Weiser, J., Wu, G., Caudy, M., Garapati, P., Gillespie, M., Kamdar, M.R., et al. (2014). The Reactome pathway knowledgebase. *Nucleic Acids Res.* **42**, D472–D477.
- DCC_portal (2016). The BLUEPRINT DCC portal. <http://dcc.blueprint-epigenome.eu>.
- DEEP (2016). Welcome to DEEP. <http://www.deutsches-epigenom-programm.de/>.
- ENCODE (2016). The ENCODE Project: ENCyclopedia of DNA elements. <https://www.genome.gov/encode/>.
- EPICO-data-loading-scripts (2016). EPICO database loading scripts GitHub repository. <https://github.com/inab/EPICO-data-loading-scripts>.
- EPICO-data-model (2016). EPICO data model GitHub repository, designed using BP-Schema-tools. <https://github.com/inab/EPICO-data-model>.
- EPICO-REST-API (2016). EPICO REST API GitHub repository. <https://github.com/inab/EPICO-REST-API>.
- Flicke, P., Amode, M.R., Barrell, D., Beal, K., Billis, K., Brent, S., Carvalho-Silva, D., Clapham, P., Coates, G., Fitzgerald, S., et al. (2014). Ensembl 2014. *Nucleic Acids Res.* **42**, D749–D755.
- Harrow, J., Frankish, A., Gonzalez, J.M., Tapanari, E., Diekhans, M., Kokocinski, F., Aken, B.L., Barrell, D., Zadissa, A., Searle, S., et al. (2012). GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res.* **22**, 1760–1774.
- ICGC (2016). International cancer genome consortium. <http://icgc.org>.
- ICGC-DCC-Docs (2016). ICGC DCC documents. <http://docs.icgc.org/>.
- IHEC (2016). International human epigenome consortium. <http://ihc-epigenomes.org/>.
- Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, D. (2002). The human genome browser at UCSC. *Genome Res.* **12**, 996–1006.
- Kurotaki, D., Yamamoto, M., Nishiyama, A., Uno, K., Ban, T., Ichino, M., Sasaki, H., Matsunaga, S., Yoshinari, M., Ryo, A., et al. (2014). IRF8 inhibits C/EBP α activity to restrain mononuclear phagocyte progenitors from differentiating into neutrophils. *Nat. Commun.* **5**, 4978.

- Malone, J., Holloway, E., Adamusiak, T., Kapushesky, M., Zheng, J., Kolesnikov, N., Zhukova, A., Brazma, A., and Parkinson, H. (2010). Modeling sample variables with an Experimental Factor Ontology. *Bioinformatics* 26, 1112–1118.
- Meehan, T.F., Vasilevsky, N.A., Mungall, C.J., Dougall, D.S., Haendel, M.A., Blake, J.A., and Diehl, A.D. (2013). Ontology based molecular signatures for immune cell types via gene expression analysis. *BMC Bioinformatics* 14, 263.
- Roadmap Epigenomics Project (2016). The NIH Roadmap Epigenomics Mapping Consortium. <http://www.roadmapepigenomics.org/>.
- Rodríguez, J.M., Maietta, P., Ezkurdia, I., Pietrelli, A., Wesselink, J.-J., Lopez, G., Valencia, A., and Tress, M.L. (2013). APPRIS: annotation of principal and alternative splice isoforms. *Nucleic Acids Res.* 41, D110–D117.
- Sarntivijai, S., Lin, Y., Xiang, Z., Meehan, T.F., Diehl, A.D., Vempati, U.D., Schürer, S.C., Pang, C., Malone, J., Parkinson, H., et al. (2014). CLO: The cell line ontology. *J. Biomed. Semantics* 5, 37.
- Smith, B., Ashburner, M., Rosse, C., Bard, J., Bug, W., Ceusters, W., Goldberg, L.J., Eilbeck, K., Ireland, A., Mungall, C.J., et al.; OBI Consortium (2007). The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration. *Nat. Biotechnol.* 25, 1251–1255.
- Ye, R.D., Boulay, F., Wang, J.M., Dahlgren, C., Gerard, C., Parmentier, M., Serhan, C.N., and Murphy, P.M. (2009). International Union of Basic and Clinical Pharmacology. LXXIII. Nomenclature for the formyl peptide receptor (FPR) family. *Pharmacol. Rev.* 61, 119–161.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited Data		
BLUEPRINT Data Release, August 2016	BLUEPRINT Consortium	ftp://ftp.ebi.ac.uk/pub/databases/blueprint/releases/20160816/
Homo Sapiens Genome Assembly Model report, release GCA_000001405.15, for GRCh38	NCBI Genome Assembly Model	ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/vertebrate_mammalian/Homo_sapiens/all_assembly_versions/GCA_000001405.15_GRCh38/GCA_000001405.15_GRCh38_assembly_report+ucsc_names.txt
Ensembl Human, release79	Flicek et al., 2014	ftp://ftp.ensembl.org/pub/release-79/mysql/homo_sapiens_core_79_38/
GENCODE, release 22	Harrow et al., 2012	ftp://ftp.sanger.ac.uk/pub/gencode/Gencode_human/release_22/
REACTOME, release 52	Croft et al., 2014	http://www.reactome.org/download/archive/52.tgz
APPRIS, release v8.2Apr2015	Rodriguez et al., 2013	http://apprisws.bioinfo.cnio.es/pub/data/homo_sapiens/ens79.v8.2Apr2015/
Cell Ontology, release 2015-05-12	Smith et al., 2007	https://raw.githubusercontent.com/obophenotype/cell-ontology/c02b0a6cbd2eaad2cbca3c402ec6086b8a1f7783/src/ontology/cl-basic.obo
Cell Line Ontology, commit 62	Samtivist et al., 2014	https://raw.githubusercontent.com/CLO-ontology/CLO/c181be5406276f8050e89d5915b40b7edbb6fc0b/src/ontology/clo.owl
UBERON ontology, release 2015-05-25	http://uberon.org/	https://github.com/obophenotype/uberon/raw/7229dc9f261da992c42685fb193465f9f544bd79/basic.obo
Phenotypic qualities (properties) ontology, commit 91	http://wiki.obofoundry.org/wiki/index.php/PATO:Main_Page	https://pato.googlecode.com/svn-history/r91/trunk/quality.obo
Experimental Factor Ontology, release 2.61 (commit 355)	Malone et al., 2010	http://sourceforge.net/p/efo/code/355/tree/trunk/src/efoinobo/efo.obo?format=raw
NCI Metathesaurus, release 201604	https://ncimeta.nci.nih.gov/	ftp://ftp1.nci.nih.gov/pub/cacore/EVS/NCI_Thesaurus/archive/nci_code_cui_map/nci_code_cui_map_201604.dat
ISO3166 ontology, release 3.59	ISO	3.59
Software and Algorithms		
BLUEPRINT / EPICO data loading scripts, release 20160819	This paper	https://github.com/inab/EPICO-data-loading-scripts/tree/20160819
Elasticsearch	https://www.elastic.co	1.7.5
Perl	https://www.perl.org	5.20
Perl module Archive::Zip	https://www.cpan.org	1.53
Perl module boolean	https://www.cpan.org	0.45
Perl module DateTime::Format::ISO8601	https://www.cpan.org	0.08
Perl module Digest::SHA1	https://www.cpan.org	2.13
Perl module Encode	https://www.cpan.org	2.73
Perl module File::Temp	https://www.cpan.org	0.2304
Perl module File::Which	https://www.cpan.org	1.09
Perl module Log::Log4perl	https://www.cpan.org	1.46
Perl module XML::LibXML	https://www.cpan.org	2.0121
Perl module URI	https://www.cpan.org	1.71
Perl module Config::IniFiles	https://www.cpan.org	2.88
Perl module JSON	https://www.cpan.org	2.90

(Continued on next page)

Continued

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Perl module Sys::CPU	https://www.cpan.org	0.61
Perl module Search::Elasticsearch	https://www.cpan.org	2.02
Perl module Tie::IxHash	https://www.cpan.org	1.23
Perl module Net::FTP::AutoReconnect	https://www.cpan.org	0.3
Perl module Set::Scalar	https://www.cpan.org	1.29
Perl module SQL::Statement	https://www.cpan.org	1.405
EPICO REST API, release v1.0.0	This paper	https://github.com/inab/EPICO-REST-API/tree/v1.0.0
Perl module Dancer2	http://perldancer.org/	0.200001
Perl module Plack::Middleware::CrossOrigin	https://www.cpan.org	0.012
Perl module Plack::Middleware::Deflater	https://www.cpan.org	0.12
Perl module FCGI	https://www.cpan.org	0.77
BLUEPRINT / EPICO Data Analysis Portal, release v1.0.4	This paper	https://github.com/inab/epico-data-analysis-portal/tree/v1.0.4
Ruby	https://www.ruby-lang.org	2.0
Ruby module compass	http://compass-style.org/	1.0.3
NodeJS	https://nodejs.org/	4.6.0
Bower package es5-shim	https://bower.io/	4.5.7
Bower package json3	https://bower.io/	3.3.2
Bower package angular	https://angularjs.org/	1.5.4
Bower package d3	https://d3js.org/	3.5.16
Bower package bootstrap-sass-official	https://bower.io/	3.3.6
Bower package angular-resource	https://bower.io/	1.5.4
Bower package angular-cookies	https://bower.io/	1.5.4
Bower package angular-sanitize	https://bower.io/	1.5.4
Bower package angular-animate	https://bower.io/	1.5.4
Bower package angular-touch	https://bower.io/	1.5.4
Bower package angular-route	https://bower.io/	1.5.4
Bower package angular-bootstrap	https://angular-ui.github.io/bootstrap/	1.3.3
Bower package select2-bootstrap-css	https://bower.io/	1.4.6
Bower package angular-ui-select	https://bower.io/	0.16.1
Bower package angular-tree-control	https://github.com/jmfernandez/angular-tree-control	6c71da1d8fb8b021af1e900cc53483b08ef3bb55
Bower package highcharts-release	http://www.highcharts.com/	v4.2.7
Bower package highcharts-ng	https://github.com/inab/highcharts-ng	5a3d1168caa5b47031f57d920d0752859bb0dc95
Bower package highcharts-export-csv	https://github.com/highcharts/export-csv	7e0b0515a52519cf14908e476c5914650b50ae38
Bower package jalette	https://github.com/emersion/jalette	b9b4108c652baef9d77ef4f32d8f09000c7dc2d8
Bower package simple-statistics	http://simplestatistics.org/	1.0.0
Bower package angular-plotly	https://bower.io/	0.1.3
Bower package roboto-fontface	https://bower.io/	0.4.5
Bower package ng-csv	http://ngmodules.org/modules/ng-csv	0.3.6
Other		
Website for the BLUEPRINT Data Analysis Portal	This paper	http://blueprint-data.bsc.es/
Installation procedures of BLUEPRINT Data Analysis Portal	This paper	https://github.com/inab/epico-data-analysis-portal/wiki

CONTACT FOR REAGENT AND RESOURCE SHARING

Alfonso Valencia contact: valencia@cniio.es

EXPERIMENTAL MODEL AND SUBJECT DETAILS

The experiments, datasets, and primary analysis that support the BLUEPRINT Data Analysis Portal are available at <http://dcc.blueprint-epigenome.eu/#/home> and <http://www.blueprint-epigenome.eu/>

METHOD DETAILS

EPICO Cyber-Infrastructure Description

The storage and querying of epigenomic data poses significant challenges to analysis portals. The main problem is that data are semi-structured and not fully organized. The analysis pipelines uses as input the results obtained by the primary analysis done by the consortium from large epigenomic experiments, ChIP-Seq histone peaks, consolidated methylated regions and gene/transcript expression values. All this information is difficult to index with traditional database due to both the size and nature of the datasets.

EPICO requires a file index with the metadata, including donor, specimen (blood or other tissue types), sample identifiers, status (healthy or disease), cellular type (referred to specimen i.e., neutrophil, monocyte, etc) and the paths to access the files with the results from the analyses. The EGA/IEC XML files by sample and experiment are required, with a description of the sample identifier, information about origin, a description about the experiment type (i.e., ChIP-Seq, WGBS or RNA-Seq) and the type of analysis performed.

EPICO data model includes the necessary sample tracking metadata (donors, specimens and samples), along with the details of the experiments performed (i.e., chromatin accessibility, WGBS, MeDIP-Seq, ChIP-Seq, mRNA-Seq and others). The results of the primary analysis pipelines have their analysis identifiers (IDs), their corresponding metrics (z-score, $-\log_{10}$ q-value, FPKM and methylation levels) and their genomic locations. These results may be a genomic region, the EnsEMBL gene ID or the EnsEMBL transcript ID with its associated metrics. Each result is mapped to its physical genomic region and linked to the corresponding metadata, such that EPICO web client can follow the path from the consolidated methylated regions, regulatory regions, expression, or histone peaks to the samples or donors through the analysis and experiments.

EPICO platform relies on a NoSQL database infrastructure to handle large volumes of semi-structured data to be stored. The EPICO data model validation is a key step in cases of unstructured data usually associated to the insertions in databases (i.e., requiring strict types, range values restriction, check valid values against a controlled vocabulary among others). We have developed the EPICO infrastructure, which take into account both the EPICO conceptual model and the physical database technology, applying concepts from object oriented programming and extended entity relationship (EER) model (Chen, 1976; Codd, 1979). The bridge data model describes the concepts, specifications and restrictions that must be validated before storing the results from the analysis pipelines. Both the results and metadata are stored in a NoSQL database instance according to the definitions and restrictions of the data model (for instance, controlled vocabularies and ontologies). EPICO software components are open access, and they are available at <https://github.com/inab>

The EPICO platform is comprised of the following modules:

- The data model (EPICO-data-model, 2016), which was initially inspired on earlier data models from ICGC DCC (ICGC-DCC-Docs, 2016), and the validation logic (BP-Schema-tools, 2015).
- The data validation and loading programs (EPICO-data-loading-scripts, 2016), which have specific parts for each project (e.g., BLUEPRINT). These specific parts store the public data results produced by each project for genes and transcripts. The generic programs also store the mappings of genes to complexes, reactions and metabolic pathways registered on REACTOME (Croft et al., 2014), as well as additional public data (for instance, principal isoform, start and stop codons or TSSs) from EnsEMBL (Flicek et al., 2014), GENCODE (Harrow et al., 2012) and APPRIS (Rodriguez et al., 2013).
- The underlying database, where all the entries are kept, as well as a copy of the data model and the controlled vocabularies (used by EPICO web client). The NoSQL database system used for BLUEPRINT deployments is Elasticsearch.
- The EPICO REST API, which implements the database queries, providing a programmatic access to the data independent from the database technology (EPICO-REST-API, 2016)
- The data analysis portal itself (BP-Data_analysis, 2016), which fetches the data relevant to the queries, consolidating them on the fly in mean data series for the different charts and views. EPICO web client has been built using web technologies and libraries, including HTML5, SVG, ES5; AngularJS, UI Bootstrap, Jallite, Ng-CSV, D3, Highcharts, Simple Statistics, Plotly and Angular Plotly.

BLUEPRINT Data Analysis Portal General Usage Description

In the first step (Figures S1A and S3A) the user introduces the query of interest, which can be a combination of genomic coordinate ranges, gene names, transcripts, exons, TSS, complexes and pathway identifiers or names from EnsEMBL (Flicek et al., 2014),

GENCODE (Harrow et al., 2012) or REACTOME (Croft et al., 2014). In addition, the user can define a flanking window around the feature to explore the genomic context. In the case of complex features, like pathways or complexes, the genomic regions of all the involved genes are shown. Moreover, the search box understands a basic search language that includes coupling (e.g., “gene:BRCA1 + gene:BRCA2”) and difference operators (e.g., “pathway:Cell Cycle, Mitotic – gene:PLK1”).

In the second step BDAP shows a set of tabs with all the genomic regions obtained from the initial user query. In each one of these tabs the web platform shows the query, a textual description of the genome layout, (i.e., the genes and transcripts in that region), and it allows the metadata of the samples and the products of the analysis related to the query to be inspected and saved. The 62 primary cell types from the hematopoietic lineage tree involved into 2016-08 release are mapped to Cell Ontology terms (Smith et al., 2007), as implemented in EPICO. The cell lines are described in EFO (Experimental Factor Ontology; Malone et al., 2010) and Cell Line Ontology (Sarntivijai et al., 2014), and the terms have links to the corresponding ontology descriptions. These ontologies do not necessarily reflect all the aspects of differentiation, but as they are built connecting each term (cell type) by relationships of “is_a” and “develops_from” (Bard et al., 2005), the ontology structures in part recapitulate the hematological differentiation. For instance, Cell Ontology has been used to study cell identity along the hematological differentiation pathways (Meehan et al., 2013).

The user selects the cell types and cell lines of his/her interest from the simplified ontology trees (Figures S1B and S3B). The number of samples is shown in function of the cell type and name for each of the selected terms, along with the results of the query region. Procedures that facilitate easy ontological sub-tree selection and ancestor de-selection are implemented in the system.

Finally, in the third step the user can select the initial analytical charts corresponding to the cell types for which epigenomics information is available, including ChIP-Seq (H3K4me3, H3K4me1, H3K27ac, H3K36me3, H3K27me3, H3K9me3 or H2A.Zac), DNase-Seq, WGBS (hypo and hyper-methylated regions) and RNA-Seq (gene or transcript level; Figures S1C and S3C). Also, the user can choose to focus on a specific disease in the shown cell types, as long as there are samples from those cell types diagnosed by that disease.

Once the three-step process is completed BDAP displays the results in a graphic form and in tables (BP-analysis, 2016). Results are distributed in three different views, the “General View” uses the different cell types as the data series for the plots, and allows to filter by disease status. The “By tissue” view allows finer grain analysis of specific tissues of origin. The “Diseases by cellular type” view uses as data series the different diseases and normal states, and allows the comparison with a set of cellular types.

RNA-Seq data is represented in box-plots as an expression value for the whole gene or for each transcript (Fragments Per Kilobase Million; BP-FPKM, 2016) and by heatmaps used to show pairwise t test comparisons (at p value level) between pairs of cellular types on the same gene or transcript. The graphs of the ChIP-Seq, DNase-Seq and Methylation data are spline-ribbon-based scatter-plots that represent the genomic coordinates on the x axis and the averaged (solid line), minimum and maximum (shadow area) $-\log_{10}(p \text{ value})$, z-values or methylation levels for the ChIP-Seq, DNA-Seq and WGBS experiments on the y axis, respectively.

These graphs also include a graphical representation of the genomic layout on the inspected genomic coordinates as special series. Genes in the genomic region are included in this graphical representation, showing in its condensed mode the transcript, exons, UTR, CDS and TSS (start and stop codons) corresponding to the principal isoform of each gene, as identified by APPRIS (Rodríguez et al., 2013).

The graphs are distributed in a fluid grid, which adapts to different screen sizes and resolutions for a better user experience when using BDAP in mobile devices. Moreover, these graphs have interactive features, like the capacity to zoom in and out on the data (all the graphs are updated at once), show and hide the legend, and switch the shown genomic layout between the condensed representation and a complete one, which includes all the transcripts. High quality renderings of each plot can be saved using their context menu for publication in png, jpeg, pdf and svg formats. The entire data series can also be downloaded in csv and xls formats for further analysis. The disease filtering menu, the list of available charts and the list of cell types with data in the query region are on the left of the grid. The user can show or hide the data series related to each cell type by clicking on it, as well as inspect the sample names and the number of samples with data in the query region.

The user can also inspect the first data entries used to compute the visible series on a specific chart. This supporting data contains the coordinates of the ChIP-Seq and the DNase-Seq peaks, the hypo- and hypermethylated regions and/or the RNA-Seq expression levels. From this view, all the supporting data can be downloaded in a tabular format.

BDAP's Browser URL is rewritten on each query, cell type and chart selection allowing to bookmark them for later inspection.

Additional links to the BLUEPRINT DCC portal (DCC_portal, 2016) and the main web page of the project are also provided. Through these links the user can browse the raw and processed data produced by the consortium, as well as a description of the methods and results, the groups participating, and the publications associated to the BLUEPRINT data.

Step-by-Step Example of BDAP Usage: FPR1

In the first step the user introduces the query of interest, FPR1, in the search box and in this case, we selected 500 bp in the flanking window size box to also explore the gene's local upstream and downstream regions (Figure S1A). In the second step the 62 primary cell types from the hematopoietic lineage tree involved into 2016-08 release are mapped to Cell Ontology terms (Smith et al., 2007). We selected the neutrophil terms based on the cell ontology hierarchy: neutrophilic myelocyte; neutrophilic metamyelocyte; band form neutrophil; segmented neutrophil of bone marrow; and mature neutrophil (Figure S1B). Finally, in the third step, from the epigenomic information available, we selected the gene expression and pairwise t test comparisons charts (from RNA-Seq) and all the histone peaks (from ChIP-Seq experiments) for H3K27Ac, H3K36me3 and H3K4me3 filtering by normal cell types to explore the FPR1 data from the neutrophil differentiation lineage (Figure S1C).

BDAP displayed the data for *FPR1*, and for its paralogs *FPR2* and *FPR3*, as these genes overlap completely or partially with the longest *FPR1* transcript annotated in Ensembl (Figure S2).

Step-by-Step Example of BDAP Usage: IRF8

We started the search introducing the gene symbol “IRF8” in the first step and with a flanking window size of 500 bps (Figure S3A). We then selected the neutrophilic myelocyte term, the classical monocyte and the macrophage terms in the second step (Figure S3B), and for the experimental results in the third step, we selected gene and transcript expression and the histone modification charts for H3K27ac, H3K27me3 and H3K9me3 filtering by normal cell types (Figure S3C).

QUANTIFICATION AND STATISTICAL ANALYSIS

BLUEPRINT Data Analysis Portal display information from the official pipelines to analyze the data produced by the project. This information is public available across BLUEPRINT-DCC or ftp site.

Detailed information about the parameters and statistics produced by BLUEPRINT and displayed by BDAP is available at:

For hypo and hyper-methylated regions: http://dcc.blueprint-epigenome.eu/#/md/bs_seq_grch38

For ChIP-seq experiments: http://dcc.blueprint-epigenome.eu/#/md/chip_seq_grch38

For DNase-seq experiments: http://dcc.blueprint-epigenome.eu/#/md/dnase_seq_grch38

For RNA-seq experiments: http://dcc.blueprint-epigenome.eu/#/md/rna_seq_grch38

The statistics displayed in the boxplot charts (median, mean, quartiles, maximum, minimum and outliers) are computed and represented with Plotly library (<https://plot.ly/>)

The Welch's t test calculated for the pairwise t test comparison charts is done by Simple Statistics library (<http://simplestatistics.org/>). These charts present the p value resulted from the test without cut-off selection. The t test comparisons are made on the fly.

DATA AND SOFTWARE AVAILABILITY

The BLUEPRINT Data Analysis Portal is available at <http://blueprint-data.bsc.es/#/>

The EPICO components can be downloaded at:

1. EPICO-data-model: <https://github.com/inab/EPICO-data-model>
2. EPICO-data-loading-scripts: <https://github.com/inab/EPICO-data-loading-scripts>
3. EPICO-REST-API: <https://github.com/inab/EPICO-REST-API>

The BLUEPRINT Data Analysis Portal components can be downloaded at:

1. BP-Schema-Tools: <https://github.com/inab/BP-Schema-tools>
2. Epico-data-analysis-portal: <https://github.com/inab/epico-data-analysis-portal>

ADDITIONAL RESOURCES

The EPICO guide installation and usage is available at <https://github.com/inab/epico-data-analysis-portal/wiki>

A first steps tutorial about BLUEPRINT Data Analysis Portal usage is available at <http://blueprint-data.bsc.es/#/first-steps>

An example usage of BLUEPRINT Data Analysis Portal is available at: <http://blueprint-data.bsc.es/#/doing-a-search>

The experiments, datasets and primary analysis that support the BLUEPRINT Data Analysis Portal are available at <http://dcc.blueprint-epigenome.eu/#/home>

The BLUEPRINT consortium data portal is available at <http://www.blueprint-epigenome.eu/>